

University of Wollongong

## Research Online

---

Faculty of Engineering and Information  
Sciences - Papers: Part B

Faculty of Engineering and Information  
Sciences

---

2019

# The Vulnerability of Multiplicative Noise Protection to Correlation-Attacks on Continuous Microdata

Yue Ma

*University of Wollongong, ym894@uowmail.edu.au*

Yan-Xia Lin

*University of Wollongong, yanxia@uow.edu.au*

Rathin Sarathy

*Oklahoma State University*

Follow this and additional works at: <https://ro.uow.edu.au/eispapers1>



Part of the [Engineering Commons](#), and the [Science and Technology Studies Commons](#)

---

### Recommended Citation

Ma, Yue; Lin, Yan-Xia; and Sarathy, Rathin, "The Vulnerability of Multiplicative Noise Protection to Correlation-Attacks on Continuous Microdata" (2019). *Faculty of Engineering and Information Sciences - Papers: Part B*. 2907.

<https://ro.uow.edu.au/eispapers1/2907>

Research Online is the open access institutional repository for the University of Wollongong. For further information contact the UOW Library: [research-pubs@uow.edu.au](mailto:research-pubs@uow.edu.au)

---

# The Vulnerability of Multiplicative Noise Protection to Correlation-Attacks on Continuous Microdata

## Abstract

When multiplicative noises are used to protect values of a sensitive attribute in a microdata, it is frequently assumed that data intruders use the noise-multiplied value to estimate the corresponding unobservable original value of a target record. In this paper, we show that, data intruders could easily construct another estimate instead of using the noise-multiplied value to attack an original value. The new estimate, namely "correlation-attack" estimate, is obtained by exploiting the potentially high correlation between the noise-multiplied data and the original data. We provide a detailed comparison between the two estimates (noise-multiplied value and the correlation-attack estimate) by comparing the mean squared errors of the two underlying estimators, and we propose that data providers should always assess the disclosure risks from both estimators when generating noise-multiplied data. Correspondingly, we propose a disclosure risk measure which could be used by data providers for noise generating variable selection during data masking stage. A simulation study is provided to illustrate how the disclosure risk measure could help with noise generating variable selection for masking a set of original data.

## Disciplines

Engineering | Science and Technology Studies

## Publication Details

Ma, Y., Lin, Y. & Sarathy, R. (2019). The Vulnerability of Multiplicative Noise Protection to Correlation-Attacks on Continuous Microdata. *Sankhya B*, Online First 1-23.

# The Vulnerability of Multiplicative Noise Protection to Correlation-Attacks on Continuous Microdata

Yue Ma · Yan-Xia Lin · Rathin Sarathy

Received: date / Accepted: date

**Abstract** When multiplicative noises are used to protect a sensitive attribute of records in a microdata, it is frequently assumed that data intruders use the noise-multiplied value to estimate the corresponding unobservable original value of a target record. In this paper, we show that, data intruders could use another estimate instead of using the noise-multiplied value to attack an original value. The new estimate, namely correlation-attack estimate, is obtained by exploiting the potentially high correlation between the noise-multiplied data and the original data. We provide a detailed comparison between the two estimates (noise-multiplied value and correlation-attack estimate) by comparing the mean square errors of the two underlying estimators, and we propose that data providers should always assess the disclosure risks from both estimators when generating noise-multiplied data. Correspondingly, we propose a disclosure risk measure which could be used by data providers for noise generating variable selection during data masking stage. A simulation study is provided to illustrate how the disclosure risk measure could be used.

**Keywords** Data Confidentiality · Noise Multiplication Masking · Continuous Microdata · Disclosure Risk · Attacking Strategy

## 1 Introduction

Microdata contains details of individuals or businesses across several attributes, such as personal income. The role of a data provider may involve collecting and releasing microdata to data users for analysis. When carrying out this role, the data provider needs to ensure that the released data carries enough information about the population, while private information of the survey respondents are not revealed to the public. To satisfy both requirements, the data provider may

---

Yue Ma  
National Institute for Applied Statistics Research, Schoole of Mathematics and Applied Statistics, University of Wollongong, Australia  
E-mail: ym894@uowmail.edu.au

Yan-Xia Lin  
National Institute for Applied Statistics Research, Schoole of Mathematics and Applied Statistics, University of Wollongong, Australia  
E-mail: yanxia@uow.edu.au

Rathin Sarathy  
Spears School of Business, Oklahoma State University, Stillwater OK 74078, USA  
E-mail: rathin.sarathy@okstate.edu

apply data perturbation techniques to produce a set of perturbed microdata, and release the perturbed microdata to the public. This paper considers the case of using multiplicative noises to perturb the original data, and the resulted noise-multiplied data is released to the public for analysis (Hwang 1986; Evans 1996; Kim and Winkler 2003; Nayak et al. 2011; Sinha et al. 2011).

Using multiplicative noises to perturb sensitive values has been advocated by many researchers because of its appealing features. Multiplicative noises provide uniform protections, in terms of the coefficient of variation of the noises, to all sensitive observations (Nayak et al. 2011) and are more suitable for economic modeling of income data (Kim and Winkler 2003). The masking mechanism is easy to implement in practice and a balanced utility-risk tradeoff is achieved by selecting an appropriate noise generating variable, which is referred to as “tuning mechanism” in Klein et al. (2014). The masking method has also been used in practice by the U.S. Energy Information Administration and the U.S. Bureau of Census (Kim and Jeong 2008).

When multiplicative noises are used to perturb a sensitive attribute, the data provider uses a noise generating variable  $C$  to produce noise terms. The data provider releases the noise-multiplied data together with some extra information of  $C$  to the public. Methodologies for analysing noise-multiplied data have been developed by taking into account the information of  $C$  so that statistical properties of the population could be recovered from the noise-multiplied microdata (Kim and Winkler 2003; Nayak et al. 2011; Sinha et al. 2011; Lin and Wise 2012; Klein et al. 2014).

It is important for data providers to understand the disclosure risk of releasing noise-multiplied microdata to the public. Disclosure occurs if the private information of a target record is learned by a data intruder. For a numeric private value (such as personal income), even if it is protected by a multiplicative noise, disclosure could still occur if the original value is reasonably inferred by a data intruder. This type of disclosure is called value disclosure or predictive disclosure (Nayak et al. 2011). Understanding all potential disclosure risks is important as it could help data providers to take precautions when generating perturbed data. However, as noted in Nayak et al. (2011), evaluating value disclosure risk is difficult because different data intruders may have different target values as well as different prior knowledge about the original data. To understand potential disclosure risks of a data masking mechanism, a common approach in the literature is to model intrusion behaviors (Agrawal and Srikant 2000; Domingo-Ferrer et al. 2004; Liu et al. 2008; Ma et al. 2016). Correspondingly, appropriate actions could be made during data masking stage such that the released data is protected against these behaviours.

Under the context of multiplicative noise perturbation, some intrusion behaviours have also been modelled under specific contexts. For instance, Nayak et al. (2011) considered the scenario that multiplicative noises are used to perturb contributor values of a tabular data, and modelled an intrusion behaviour for disclosing a target contributor value of a table cell. In terms of using multiplicative noises to perturb microdata, Klein et al. (2014) considered the scenario where a response variable is sensitive and masked by multiplicative noises while explanatory variables are unmasked. The authors showed that a data intruder may use the predicted value based on the generalised regression model to estimate a private original value. Another intrusion behaviour is recognised in Nayak et al. (2011) and Lin and Wise (2012), that the data provider may simply use the noise-multiplied value as an estimate of an original value because the noise-multiplied value is in fact an unbiased estimate of the original value. This fact is also recognised in Kim (2007) and Kim and Jeong (2008) and these authors proposed several noise distributions in order to reduce the disclosure risk. Corresponding disclosure risk measures against these intrusion behaviours were proposed in Klein et al. (2014) and Lin and Wise (2012) so that data providers could control the disclosure risks by choosing a suitable noise candidate during data masking stage.

This paper considers another intrusion behaviour for attacking noise-multiplied data. Similar to the intrusion behaviour recognised in Nayak et al. (2011) and Lin and Wise (2012), the

intrusion behaviour we consider in this paper only uses information of the noise-multiplied attribute itself. Because of its simplicity, the attacking strategy could always be used as an attempt to unveil a noise-multiplied value. The attacking strategy exploits the correlation between the original data and the noise-multiplied data, therefore we name it as “correlation-attack”. The intuition behind the correlation-attack is that if the correlation between the original data and the noise-multiplied data is high, then a simple linear regression model might be adequate to explain the relationship between the two sets of data. The correlation between the original data and the noise-multiplied data could be reasonably estimated by a data intruder, as we show in Section 3. On the basis of the mean squared errors (MSEs), both the unbiased estimator and the correlation-attack estimator tend to be more accurate as the correlation gets larger. Therefore, a high correlation value might motivate a data intruder to use either the unbiased estimate or the correlation-attack estimate to disclose a target value. Correspondingly, we propose a disclosure risk measure to be used by data providers, so that for each original value, the disclosure risks from both estimators could be simultaneously evaluated. The disclosure risk measure could help data providers with noise generating variable selection during data masking stage. A simulation study is provided to illustrate the use of the disclosure risk measure.

The idea of using regression models to attack private values has also been considered for k-anonymized data (Li and Sarkar 2011). The authors considered the possibility of using a regression tree to attack a protected value of a target record. By treating the target anonymized attribute as the response variable and other anonymized attributes as explanatory variables, the authors showed that it is possible to reveal the private information of a target record. The authors referred to this attack as “regression-attack”. The regression-attack is similar to the idea of using the predicted value based on a generalised linear regression model to attack a noise-multiplied value proposed in Klein et al. (2014). We feel that the major difference between the regression-attack and our correlation-attack is that the correlation-attack only uses information of the noise-multiplied attribute itself to attack a noise-multiplied value. That is, the regressors of a regression-attack could involve other confidentialised attributes in a microdata, while the regressor of a correlation-attack is simply the confidentialised attribute itself, or the noise-multiplied attribute in our context. We note that when several attributes in a microdata are protected by multiplicative noises, the idea of the regression-attack might be used to attack a particular noise-multiplied value. In that case, the regressors of an attacking regression model are some or all attributes (either noise-multiplied or not) in the microdata. This idea could be explored in the future.

Throughout this paper, we assume all original data and noise candidates are positive and continuous, and all noise candidates have expectation 1. The sample size of the original data is large enough so that estimates of some population parameters could be accurately recovered. We assume the following simple scenario: the data provider releases the noise-multiplied data together with the variance of the noise generating variable used to perturb a set of original data to the public, so that data users could unbiasedly recover the first two moments estimates of the population (such as population mean and variance) from the noise-multiplied data. We note that assuming preservation of the first two moments estimates only is consistent with many other masking schemes in the literature (Brand 2002; Yancey et al. 2002; Kim and Winkler 2003; Domingo-Ferrer et al. 2004; Oganian and Karr 2011). For instance, the data perturbation method proposed in Oganian and Karr (2011) only allows data users to obtain unbiased estimates of population means and covariance matrix from perturbed microdata. The assumption simplifies discussions of this paper especially when we define overall data utility loss. Sophisticated methodologies for analysing noise-multiplied data require the density functions of noise generating variables to be public (Sinha et al. 2011; Klein et al. 2014; Lin 2014; Lin and Fielding 2015), so that more population information could be recovered from noise-multiplied data. The

correlation-attack proposed in this paper could also be used under these situations where density functions of noises are public.

The paper is organized as follows: Section 2 introduces the mathematical setup of the noise multiplication masking method and review methodologies for recovering the first two moments estimates of a population from noise-multiplied data. Section 3 introduces the correlation-attack strategy. Section 4 discusses and compares the correlation-attack estimator and the unbiased estimator. Section 5 proposes a disclosure risk measure which could be used by data providers for noise candidates selection and introduces the definition of overall data utility loss we adopt in this paper. Section 6 presents a simulation study. Section 7 concludes the paper.

## 2 Recovering the first two moments population estimates from noise-multiplied data

In this section we describe the mathematical setup of noise multiplication masking method and review methodologies proposed in Nayak et al. (2011) for recovering the first two moments estimates (first two moments and variance) of a population from a set of noise-multiplied data. For a set of univariate original data, the noise multiplication masking method works as follows:

*Suppose a set of original data  $y = \{y_i\}_{i=1}^n$  are independent realizations from  $Y$ . To mask  $y$ , the data provider chooses a noise generating variable  $C$  with  $E(C) = 1$ .  $Y$  and  $C$  are independent. A set of noise terms  $c = \{c_i\}_{i=1}^n$  are independently drawn from  $C$  and multiplied with  $y$  to produce the noise-multiplied data  $y^* = \{y_i^*\}_{i=1}^n = \{y_i c_i\}_{i=1}^n$ .  $\{y_i^*\}_{i=1}^n$  could be treated as independent realizations from  $Y^*$ , where  $Y^* = YC$ . The data provider releases  $y^*$  together with other information of  $C$  (such as variance  $\sigma_C^2$  or density function  $f_C$ ) to the public. Following Kim and Winkler (2003) and Lin and Wise (2012), throughout this paper we assume a simple case where only  $\sigma_C^2$  is released to the public.*

If  $\sigma_C^2$  is public, Nayak et al. (2011) showed that the first two moments of  $Y$  could be unbiasedly recovered from  $Y^*$ . That is

$$E(Y) = E(Y^*)$$

and

$$E(Y^2) = \frac{E[(Y^*)^2]}{E(C^2)} = \frac{E[(Y^*)^2]}{(\sigma_C^2 + 1)}.$$

Denote  $E(Y) = \mu_Y$  and  $Var(Y) = \sigma_Y^2$ . The authors also showed that,  $\mu_Y$  and  $\sigma_Y^2$  could be unbiasedly estimated by  $\hat{\mu}_Y$  and  $\hat{\sigma}_Y^2$  using  $y^*$ , where

$$\hat{\mu}_Y = \frac{\sum_{i=1}^n y_i^*}{n} \quad (1)$$

and

$$\hat{\sigma}_Y^2 = \frac{1}{n(n-1)} \left[ \left( \frac{n + \sigma_C^2}{1 + \sigma_C^2} \right) \sum_{i=1}^n (y_i^*)^2 - \left( \sum_{i=1}^n y_i^* \right)^2 \right]. \quad (2)$$

In summary, when  $\sigma_C^2$  is public, data users could unbiasedly recover estimates of the first two moments and variance of  $Y$ . However, it is not true for other information of  $Y$ , such as quantiles and higher order moments as recovering these information requires  $f_C$  to be public.

### 3 The correlation-attack strategy

In this section we introduce the correlation-attack strategy which could be used by a data intruder to obtain an estimate of an unobservable original value  $y_i$ . The correlation-attack follows from the idea that, if the sample correlation between  $y^*$  and  $y$  is high, then a simple linear regression model may adequately explain the relationship between  $y$  and  $y^*$ . Consequently, the data intruder could use the predicted value based on the simple linear regression model to estimate a target original value  $y_i$  from  $y_i^*$ .

From the data intruder's perspective, the original data  $y$  is unobservable. Therefore, the sample correlation between  $y$  and  $y^*$ , denoted as  $r_{yy^*}$ , cannot be known. However, when the sample size  $n$  is large,  $r_{yy^*}$  is close to  $\rho_{YY^*}$ , where  $\rho_{YY^*}$  is the population correlation coefficient between  $Y$  and  $Y^*$ . It can be shown that  $\rho_{YY^*}$  takes the following form:

$$\rho_{YY^*} = \frac{Cov(Y^*, Y)}{\sqrt{Var(Y^*)Var(Y)}} = \sqrt{\frac{Var(Y)}{Var(Y^*)}} = \sqrt{\frac{\sigma_Y^2}{\sigma_Y^2(\sigma_C^2 + 1) + \mu_Y^2\sigma_C^2}} \quad (3)$$

In the above expression, the knowledge of  $\mu_Y$  and  $\sigma_Y^2$  is unknown. However, the data intruder could unbiasedly estimate these two terms from  $y^*$  by  $\hat{\mu}_Y$  and  $\hat{\sigma}_Y^2$  using Equation (1) and (2). Therefore,  $\rho_{YY^*}$  could be approximated by  $\tilde{r}_{yy^*}$ , where

$$\tilde{r}_{yy^*} = \sqrt{\frac{\hat{\sigma}_Y^2}{\hat{\sigma}_Y^2(\sigma_C^2 + 1) + \hat{\mu}_Y^2\sigma_C^2}}. \quad (4)$$

As a result, if the sample size  $n$  is large, then  $\tilde{r}_{yy^*}$  could be used to approximate the sample correlation  $r_{yy^*}$ . If  $\tilde{r}_{yy^*}$  is large, then it might motivate the data intruder to fit a simple linear regression model between  $y$  and  $y^*$ , and to use the predicted value based on the linear model to attack the unobservable  $y_i$ . The predicted value of  $y_i$  based on  $y_i^*$  has the following expression:

$$\hat{y}_i = \hat{\alpha} + \hat{\beta}y_i^*,$$

where  $\hat{\alpha}$  and  $\hat{\beta}$  are least squares estimates of the intercept and slope terms. The estimates can be evaluated as follows:  $\hat{\beta} = r_{yy^*} \frac{s_y}{s_{y^*}}$  and  $\hat{\alpha} = \bar{y} - \hat{\beta}\bar{y^*}$ , where  $s_y$  and  $s_{y^*}$  are sample variances of  $y$  and  $y^*$ , respectively;  $\bar{y}$  and  $\bar{y^*}$  are sample means of  $y$  and  $y^*$ , respectively.

For the data intruder,  $\bar{y}$ ,  $s_y$  and  $r_{yy^*}$  cannot be known directly as  $y$  is not available. Therefore  $\hat{\alpha}$  and  $\hat{\beta}$  cannot be obtained directly. However, when the sample size  $n$  is large,  $\bar{y}$  is close to  $\mu_Y$ ,  $s_y$  is close to  $\sigma_Y$ , and  $r_{yy^*}$  is close to  $\rho_{YY^*}$ . The values  $(\mu_Y, \sigma_Y, \rho_{YY^*})$  could be approximated from  $y^*$  by  $(\hat{\mu}_Y, \hat{\sigma}_Y, \tilde{r}_{yy^*})$  using Equation (1), (2) and (4). Therefore, the unknown quantities  $\bar{y}$ ,  $s_y$  and  $r_{yy^*}$  could be approximated by  $\hat{\mu}_Y$ ,  $\hat{\sigma}_Y$  and  $\tilde{r}_{yy^*}$  respectively. As a result, the least squares estimates  $\hat{\alpha}$  and  $\hat{\beta}$  could be approximated by  $\tilde{\alpha}$  and  $\tilde{\beta}$ , where  $\tilde{\beta} = \tilde{r}_{yy^*} \frac{\hat{\sigma}_Y}{s_{y^*}}$  and  $\tilde{\alpha} = (1 - \tilde{r}_{yy^*} \frac{\hat{\sigma}_Y}{s_{y^*}})\bar{y^*}$ . Therefore, the correlation-attack estimate  $\tilde{y}_i$  which could be obtained by the data intruder to attack  $y_i$  takes the following form:

$$\tilde{y}_i = \tilde{\alpha} + \tilde{\beta}y_i^* = (1 - \tilde{r}_{yy^*} \frac{\hat{\sigma}_Y}{s_{y^*}})\bar{y^*} + \tilde{r}_{yy^*} \frac{\hat{\sigma}_Y}{s_{y^*}}y_i^*.$$

### 4 Comparison between $Y_i^*$ and $\tilde{Y}_i$ and discussion

In this section we discuss and compare the accuracies of the correlation-attack estimator  $\tilde{Y}_i$  and the unbiased estimator  $Y_i^*$  for estimating a target value  $y_i$ . We compare the two estimators

because both  $\tilde{y}_i$  and  $y_i^*$  could easily be obtained by a data intruder given basic knowledge of  $y^*$  and  $\sigma_C^2$  to attack  $y_i$ . Because of the generality, we propose that the disclosure risks from both estimators should always be evaluated and controlled for most noise-multiplied data. To understand both estimators, we base our discussion on the mean squared errors (MSE) of these estimators.

From the last section, we introduced the correlation-attack estimate  $\tilde{y}_i$ . We denote the correlation-attack estimator as  $\tilde{Y}_i$ . The mathematical expression of  $\tilde{Y}_i$  cannot be easily found without making a large sample assumption. By noting that when sample size is large,  $\tilde{r}_{yy^*} \frac{\hat{\sigma}_Y}{s_{y^*}}$  is close to  $\rho_{YY^*} \frac{\sigma_Y}{\sigma_{Y^*}} = \rho_{YY^*}^2$ , and  $\bar{y}^*$  is close to  $\mu_Y$ . Then,  $\tilde{Y}_i$  is approximated to be the following

$$\tilde{Y}_i \approx (1 - \rho_{YY^*}^2)\mu_Y + \rho_{YY^*}^2 Y_i^*.$$

We use this expression of  $\tilde{Y}_i$  for discussion throughout this section. On the other hand, the unbiased estimator  $Y_i^*$  is simply the noise multiplied variable and  $Y_i^* = Y_i C_i$ . It is an unbiased estimator because  $E(Y_i^* | Y_i = y_i) = y_i$ . It can be seen directly that  $\tilde{Y}_i$  is partially determined by  $Y_i^*$ . The MSE of the correlation-attack estimator  $\tilde{Y}_i$  is given as:

$$MSE(\tilde{Y}_i | Y_i = y_i) = E[(1 - \rho_{YY^*}^2)\mu_Y + \rho_{YY^*}^2 C Y_i - Y_i | Y_i = y_i]^2 = (1 - \rho_{YY^*}^2)^2 (\mu_Y - y_i)^2 + \rho_{YY^*}^4 y_i^2 \sigma_C^2.$$

The MSE of the unbiased estimator  $Y_i^*$  is given as:

$$MSE(Y_i^* | Y_i = y_i) = E(C Y_i - Y_i | Y_i = y_i)^2 = y_i^2 E(C - 1)^2 = y_i^2 \sigma_C^2.$$

**Accuracies of  $\tilde{Y}_i$  and  $Y_i^*$ :** We first note that from Equation (3),  $\sigma_C^2 = \frac{1 - \rho_{YY^*}^2}{\rho_{YY^*}^2 (1 + \mu_Y^2 / \sigma_Y^2)}$ . Because  $\rho_{YY^*} \in [0, 1]$ , therefore  $\sigma_C^2$  monotonically decreases as  $\rho_{YY^*}$  increases. The MSEs of both estimators decrease as  $\sigma_C^2$  decreases, meaning that they are more accurate for estimating  $y_i$  as  $\rho_{YY^*}$  gets larger. It is straightforward to see it for  $Y_i^*$  based on its MSE. However, it is not straightforward to see it for  $\tilde{Y}_i$ . To show this, we substitute  $\sigma_C^2$  by  $\frac{1 - \rho_{YY^*}^2}{\rho_{YY^*}^2 (1 + \mu_Y^2 / \sigma_Y^2)}$  in the expression of  $MSE(\tilde{Y}_i | Y_i = y_i)$ , so we have

$$MSE(\tilde{Y}_i | Y_i = y_i) = (k_1 - k_2)\rho_{YY^*}^4 + (k_2 - 2k_1)\rho_{YY^*}^2 + k_1,$$

where  $k_1 = (\mu_Y - y_i)^2$ ,  $k_2 = \frac{y_i^2}{1+h}$ ,  $h = \mu_Y^2 / \sigma_Y^2$ . So the MSE is a parabola in terms of  $\rho_{YY^*}^2$ . The symmetric axis is  $S = 1 + \frac{k_2}{2(k_1 - k_2)}$ . The parabola is monotonically decreasing in  $\rho_{YY^*}^2 \in [0, 1]$  in almost all cases, meaning that the correlation-attack estimator  $\tilde{Y}_i$  is more accurate as  $\rho_{YY^*}$  increases. The only exception is when  $k_2 > 2k_1$ . Under this condition, the symmetric axis  $S$  is within  $[0, 1]$  so the function is not monotone in  $[0, 1]$ . That means when  $y_i \in (\frac{2\mu_Y(1+h) - \mu_Y\sqrt{2(1+h)}}{1+2h}, \frac{2\mu_Y(1+h) + \mu_Y\sqrt{2(1+h)}}{1+2h})$ ,  $MSE(\tilde{Y}_i | Y_i = y_i)$  will increase first as  $\rho_{YY^*}^2$  goes from 0 to  $S$ , but it eventually decreases to 0 as  $\rho_{YY^*}^2$  increases from  $S$  to 1. Therefore, both estimators are more accurate as  $\rho_{YY^*}$  gets larger. A large  $\rho_{YY^*}$  value might motivate a data intruder to use either the unbiased estimator or the correlation-attack estimator to attack  $y_i$ .

**Comparison between  $\tilde{Y}_i$  and  $Y_i^*$ :** An estimator predicts the unknown value better if it yields a smaller MSE. To compare which estimator predicts  $y_i$  better, we set  $MSE(\tilde{Y}_i | Y_i = y_i) - MSE(Y_i^* | Y_i = y_i) < 0$ . Those conditions which satisfy this inequality mean that under the



conditions, the correlation-attack estimator  $\tilde{Y}_i$  predicts  $y_i$  with better accuracy. After solving the inequality, we let

$$\begin{aligned} a &= \frac{\sigma_Y^2 - \sigma_Y \sqrt{\sigma_Y^2 + \mu_Y^2}}{\mu_Y^2}, \\ b &= \frac{\sigma_Y^2 + \sigma_Y \sqrt{\sigma_Y^2 + \mu_Y^2}}{\mu_Y^2}, \\ c &= \frac{\mu_Y \rho_{YY^*}^2 (\sigma_Y^2 + \mu_Y^2) - \mu_Y \sigma_Y \rho_{YY^*} \sqrt{(1 + \rho_{YY^*}^2)(\sigma_Y^2 + \mu_Y^2)}}{\rho_{YY^*}^2 \mu_Y^2 - \sigma_Y^2}, \\ d &= \frac{\mu_Y \rho_{YY^*}^2 (\sigma_Y^2 + \mu_Y^2) + \mu_Y \sigma_Y \rho_{YY^*} \sqrt{(1 + \rho_{YY^*}^2)(\sigma_Y^2 + \mu_Y^2)}}{\rho_{YY^*}^2 \mu_Y^2 - \sigma_Y^2}, \end{aligned}$$

$e = \min(c, d)$  and  $f = \max(c, d)$ , then we have the following result:

**Result 1:** Based on MSEs, for large-sized sample, when  $\rho_{YY^*} < a$  or  $\rho_{YY^*} > b$ , an observation  $y_i$  is more vulnerable to  $\tilde{Y}_i$  if  $y_i \in (e, f)$ ; when  $\rho_{YY^*} \in (a, b)$ ,  $y_i$  is more vulnerable to  $\tilde{Y}_i$  if  $y_i < e$  or  $y_i > f$ ; when  $\rho_{YY^*}$  equals  $a$  or  $b$ ,  $y_i$  is more vulnerable to  $\tilde{Y}_i$  if  $y_i > \mu_Y/2$ .

Based on Result 1, we have the following discussions from both the data intruder and the data provider's perspectives:

*Discussion 1:* If  $\rho_{YY^*}$  is high, then the data intruder may use either the correlation-attack estimator or the unbiased estimator to attack a target value  $y_i$ . The values of  $(a, b, e, f)$  depend on three parameters  $(\mu_Y, \sigma_Y, \rho_{YY^*})$ . From the data intruder's perspective, these parameters could easily be estimated from  $y^*$  using Equations (1), (2) and (4). Therefore, the data intruder could obtain estimates  $(\hat{a}, \hat{b}, \hat{e}, \hat{f})$ , and use Result 1 to make a decision on which particular estimator should be used for attacking  $y_i$ . In order to do this, the data intruder needs to make an initial guess about the location of  $y_i$  in terms of  $(\hat{e}, \hat{f})$ . For instance, if the data intruder's estimate  $\tilde{r}_{yy^*}$  is within  $(\hat{a}, \hat{b})$ , and the data intruder has a strong belief that  $y_i$  is greater than  $\hat{f}$ , then logically speaking the data intruder would use the correlation-attack estimate to attack the value of  $y_i$ . In the simulation study in Section 6, it can be shown that original values over 26317.6 are more vulnerable to the correlation-attack estimator according to Result 1. If a noise-multiplied value is significantly greater than 26317.6, say 200000, it is unlikely that the corresponding original value is less than 26317.6. Therefore, it is very likely that the data intruder uses the correlation-attack estimator to attack this value.

Alternatively, if a data intruder has no prior knowledge about the location of the unobservable  $y_i$  and only assumes  $y_i \sim Y_i$ , one possible choice is that the data intruder may attempt to compute the expected location of  $y_i$  by  $E(Y_i|Y_i^* = y_i^*) = y_i^* E(\frac{1}{C})$ , and make a decision about which estimator to use according to this value. The exact value of  $E(\frac{1}{C})$  could only be known if the density function  $f_C$  is public. Because we only assume  $\sigma_C^2$  is public,  $E(\frac{1}{C})$  cannot be known. However, it is commonly assumed in the literature that noise generating variables are symmetric and positive. Under this assumption, the support of  $C$  will be a subset of  $[0, 2]$  and

$$E(\frac{1}{C}) = \lim_{\epsilon \rightarrow 0} \int_{\epsilon}^{2-\epsilon} \frac{1}{c} f_C(c) dc = \lim_{\epsilon \rightarrow 0} \int_{\epsilon}^1 (\frac{1}{c} + \frac{1}{2-c}) f_C(c) dc \geq 2 \int_{\epsilon}^1 f_C(c) dc \geq 1,$$

which means  $E(Y_i|Y_i^* = y_i^*) = y_i^* E(\frac{1}{C}) \geq y_i^*$ . Thus, the data intruder could have a rough idea about the expected location of  $y_i$  based on  $y_i^*$ , and make a decision about which estimator to use

according to this information and Result 1. In summary, the data intruder may use either  $Y_i^*$  or  $\tilde{Y}_i$  to attack the unobservable  $y_i$  according to Result 1.

*Discussion 2:* From the data provider’s perspective, the above result means that the correlation-attack estimator  $\tilde{Y}_i$  could yield a higher disclosure risk than the unbiased estimator  $Y_i^*$  for some  $y_i$  and vice versa. To protect all observations, the safest way is to make sure that for each original value  $y_i$ , the disclosure risks from both estimators are simultaneously below an acceptable level. In that way, all original observations could be protected regardless of the data intruder’s attacking strategy.

In the next section, we propose a disclosure risk measure which could be used by data providers to measure disclosure risks against these two estimators simultaneously. The proposed disclosure risk measure could help data providers with noise generating variable selection during data masking stage.

## 5 Disclosure risk and data utility loss measures

For noise multiplication masking method, the noise generating variable  $C$  plays the role of balancing data utility loss and disclosure risk, or “tuning mechanism” in Klein et al. (2014). That is, for a set of original data, the data provider needs to decide on an appropriate noise generating variable which achieves the required utility-risk tradeoff. In this section we introduce our definition of **disclosure risk** and propose a **disclosure risk measure** to be used by data providers for noise generating variable selection in practice. We also introduce the definition of **overall data utility loss** we use in this paper

Disclosure occurs if a data intruder successfully identified a target record in the microdata, and learned a private value of the identified record. That requires both identity disclosure and value disclosure to occur at the same time. Identity disclosure occurs if a data intruder successfully identified a record through a set of quasi-identifiers (Li and Sarkar 2013), or through record-linkage techniques (Kim and Winkler 1995; Oganyan and Karr 2011). Identity disclosure risk could be reduced by certain masking techniques, such as k-anonymity (Li and Sarkar 2013). The noise multiplication masking method could reduce identification rate caused by record-linkage techniques (Muralidhar and Domingo-Ferrer 2016). However, it may not reduce identity disclosure risk caused by quasi-identifiers, especially that some quasi-identifiers are non-numeric and cannot be protected by multiplicative noises. Therefore, given a noise-multiplied microdata, we may not know which records might be identified by data intruders. In this paper, we conservatively assume that all records are vulnerable to identity disclosure. Therefore, our focus is to reduce value disclosure risk (or predictive disclosure in Nayak et al. 2011). A low value disclosure risk means that even if a record is being identified by a data intruder, the private information associated with the record could not be confidently learned by the data intruder because it is protected by a multiplicative noise. Value disclosure occurs if a data intruder’s estimate of a target original value is close to the real value. Following Lin and Wise (2012), Klein et al. (2014) and Ma et al. (2016), we define “disclosure” in the following way:

*To estimate a target value  $y_i$  (positive and continuous), a data intruder may use his own attacking strategy to obtain an estimate  $\hat{y}_i$ . To classify  $\hat{y}_i$  as a good estimate of  $y_i$ , it is sufficient for  $\hat{y}_i$  to be reasonably close to  $y_i$ . **Disclosure** of  $y_i$  occurs if  $|\frac{\hat{y}_i - y_i}{y_i}| < \delta$ , where  $\delta$  is called **acceptance rule** in Lin and Wise (2012) and is a small positive number. For instance, if  $\delta = 0.05$ , we say*

that  $\hat{y}_i$  discloses  $y_i$  if  $0.95y_i < \hat{y}_i < 1.05y_i$ .

Different measures for quantifying value disclosure risk of  $y_i$  have been proposed in the literature, with the assumption that the unbiased estimator  $Y_i^*$  is used for disclosing the original value. For instance, Nayak et al. (2011) introduced a confidence interval measure to quantify the disclosure risk. Specifically, let  $[a_1, b_1]$  be the shortest interval satisfying  $P(a_1 < C_i < b_1) = 1 - \alpha$ , where  $\alpha$  is a positive number. So the  $100(1 - \alpha)$  confidence interval of  $y_i$  is  $(y_i^*/b_1, y_i^*/a_1)$ , which is used to measure the disclosure risk of  $y_i$ . A similar measure is used in Agrawal and Srikant (2000) under the context of additive noise masking. Another measure proposed in Lin and Wise (2012) calculates the probability of  $y_i$  being disclosed by  $Y_i^*$ , which is given by

$$R_{LW}(Y_i^*, \delta | Y_i = y_i) = P(|\frac{Y_i^* - y_i}{y_i}| < \delta) = P(|C - 1| < \delta). \quad (5)$$

We propose to use the probabilistic disclosure risk measure because it is straight-forward and could easily be modified to measure disclosure risks from other estimators. For instance, Klein et al. (2014) proposed a similar probabilistic disclosure risk measure by replacing  $Y_i^*$  by a linear predictor based on a regression model. Following this idea, we propose a disclosure risk measure  $R_\rho(y_i, \delta | Y_i = y_i)$  to evaluate the disclosure risk of  $y_i$  being disclosed by the correlation-attack estimator  $\tilde{Y}_i$ . That is:

$$R_\rho(\tilde{Y}_i, \delta | Y_i = y_i) = P(|\frac{\tilde{Y}_i - y_i}{y_i}| < \delta). \quad (6)$$

As the exact form of  $\tilde{Y}_i$  is not straight-forward to show,  $R_\rho$  cannot be computed exactly. For large-sized sample, we propose two approaches for data providers to estimate  $R_\rho$  in practice.

**Approach 1:** The data provider may just assume that  $\tilde{Y}_i \approx (1 - \rho_{Y^*Y}^2)\mu_Y + \rho_{Y^*Y}^2 Y_i^*$  as in Section 4. Then, we have  $R_\rho(\tilde{Y}_i, \delta | Y_i = y_i) \approx P(\frac{(-\delta+1)y_i - (1-\rho_{Y^*Y}^2)\mu_Y}{\rho_{Y^*Y}^2 y_i} < C < \frac{(\delta+1)y_i - (1-\rho_{Y^*Y}^2)\mu_Y}{\rho_{Y^*Y}^2 y_i})$ . In this expression, the parameters  $\mu_Y$  and  $\rho_{Y^*Y}$  are not known to the data provider but they could be estimated using sample estimates obtained from the original data. That is,  $\mu_Y$  could be estimated by the sample mean of the original data, and the parameter  $\rho_{Y^*Y}$  could be estimated by plugging in sample estimates of  $\mu_Y$  and  $\sigma_Y^2$  calculated from the original data into Equation (??). The data provider could then use the approximated disclosure risk measure to estimate  $R_\rho$ . Our simulations showed that for large-sized sample, e.g. sample size greater than 1000, the approximates of  $R_\rho$  are very close to the theoretical values for each  $y_i$ .

**Approach 2:** The data provider could use Monte-Carlo simulations to approximate  $R_\rho$ . That is, the data provider could firstly produce  $N$  multiple copies of noise-multiplied data using the noise candidate  $C$ , and then perform the correlation-attack by assuming the role of a data intruder. The correlation-attack is applied on each copy of the noise-multiplied data following the steps described in Section 3. To estimate  $R_\rho$ , suppose among the  $N$  copies, in  $q$  of them  $y_i$  is disclosed by the corresponding correlation-attack estimate. Then  $R_\rho(\tilde{Y}_i, \delta | Y_i = y_i)$  is estimated to be  $q/N$ .

Hereafter, we use  $R_{LW}(y_i, \delta)$  to denote  $R_{LW}(Y_i^*, \delta | Y_i = y_i)$  etc. We use  $R_{cor}(y_i, \delta)$  to denote an approximate of  $R_\rho(y_i, \delta)$  regardless of which approach the data provider uses to estimate  $R_\rho(y_i, \delta)$ . From the last section, we have shown that some original observations are more vulnerable to the correlation-attack estimator while the others are more vulnerable to the unbiased estimator. We note that even though Result 1 provides a guidance about which estimator is more effective for predicting each  $y_i$ , it is calculated based on MSEs. Under the probabilistic disclosure risk measure we proposed, our simulation results with different sets of synthetic data show that

an estimator with a lower MSE does not always result in a larger disclosure risk. Therefore, to protect  $y_i$ , the safest way is to make sure that the disclosure risks from both estimators are simultaneously below an acceptable level. In that way, we say that  $y_i$  is protected regardless of the data intruder's attacking strategy. We propose the following disclosure risk measure to evaluate the disclosure risk for each  $y_i$ :

$$R(y_i, \delta) = \max\{R_{LW}(y_i, \delta), R_{cor}(y_i, \delta)\} \quad (7)$$

A sufficiently low  $R(y_i, \delta)$  value means that  $y_i$  is protected against both the correlation-attack estimator and the unbiased estimator. For a set of original data  $\{y_i\}_{i=1}^n$  and a noise candidate  $C$ , the data provider collects a set of disclosure risks  $\{R(y_i, \delta)\}_{i=1}^n$  for each  $y_i$ . The data provider could use the set of disclosure risks as a reference for noise candidate selection. In some cases, it might be sufficient to say that a noise candidate offers an acceptable level of protection to the original data if the average value  $\overline{\{R(y_i, \delta)\}_{i=1}^n}$  is below a threshold value. In some other cases where all observations are highly sensitive, the data provider may require  $\max(\{R(y_i, \delta)\}_{i=1}^n)$  to be sufficiently low in order for a noise candidate to be considered. The data provider might need to determine a criteria for disclosure risk control according to the nature of the original data. We will provide an example in the simulation study.

We note that the disclosure risk measure is designed for the data provider to understand the disclosure risk of using a noise candidate  $C$  to mask the original data. For the data intruder, he/she may rely on Discussion 1 to determine which attacking estimator to use, but the intruder cannot use the disclosure risk measure to evaluate the probability of successfully disclosing  $y_i$  using either  $Y_i^*$  or  $\tilde{Y}_i$ . We also note that a low  $R(y_i, \delta)$  value guarantees that  $y_i$  is protected against both the unbiased estimator  $Y_i^*$  and the correlation-attack estimator  $\tilde{Y}_i$ . However, we may not say that  $y_i$  has a low disclosure risk because of other possible attacking strategies. When  $y_i$  is subject to other disclosure risks, the data provider might need to consider using  $R(y_i, \delta)$  together with other disclosure risk measures to jointly determine a suitable noise generating variable. Because the unbiased estimator and the correlation-attack estimator are two basic attacking estimators, we feel that a noise candidate should at least guarantee that the noise-multiplied data is protected against these two attacking estimators first in order for the noise candidate to be considered further for masking a set of original data.

Even though the original data is better protected if  $\sigma_C^2$  is larger,  $\sigma_C^2$  cannot be too large as doing so may significantly reduce the analytical validity of the noise-multiplied data, or data utility. When recovering a population parameter estimate from the noise-multiplied data, the recovered parameter estimate is less accurate than the one data users would obtain by analysing the original data. Data utility loss measures such loss of accuracy for a population parameter. In the following we name the recovered parameter estimates from noise-multiplied data as perturbed estimates, and the estimates obtained by analysing the original data as unperturbed estimates. An **overall data utility loss** is an aggregate measure of the utility losses across several parameters. In the remaining of this section, we introduce the overall data utility loss measure we use in this paper.

There is no unique way to measure overall data utility loss. In the literature, the way of measuring overall data utility loss varies according to different data masking scenarios as well as which parameters estimates could be recovered from masked data. For instance, Yancey et al. (2002) proposed to use the average of relative distances between perturbed estimates and unperturbed estimates across several parameters as an overall data utility loss measure under the context of additive-noise perturbation. Those parameters include population means, covariances and correlations, as those parameter estimates could be accurately recovered from noise-added microdata. Other overall utility loss measures are available, see Shlomo (2010), Domingo-Ferrer and Torra (2001) and Agrawal and Aggarwal (2001).

Duncan et al. (2001; 2004) proposed to use a data user's mean squared error in estimating a population parameter from perturbed data as a way to measure utility loss. We note that, when noise-multiplied data  $y^*$  and noise variance  $\sigma_C^2$  are released to the public, data users could only recover the first two moments estimates, such as population mean and variance, from  $y^*$ . Calculating these estimates requires estimates of the first two moments of  $Y$ . In this paper, we assume a simple case that the overall utility loss measure is computed according to utility losses of  $E(Y)$  and  $E(Y^2)$  only. We note that similar cases for measuring overall utility loss have also been considered in the literature (Kim and Winkler 1995, 2003; Brand 2002; Yancey et al. 2002; Oganian and Karr 2011). For instance, Kim and Winkler (1995) only considered utility losses for population mean and variance as indications of overall utility loss for noise-added microdata.

In the following we introduce the overall utility loss measure we use in this paper. From Section 2,  $E(Y) = E(Y^*)$  and  $E(Y^2) = \frac{E[(Y^*)^2]}{(\sigma_C^2 + 1)}$ . Therefore,  $E(Y)$  could be unbiasedly estimated by data users using  $U_Y = \frac{\sum_{i=1}^n Y_i^*}{n}$ , and  $E(Y^2)$  could be unbiasedly estimated using  $U_{Y^2} = \frac{\sum_{i=1}^n (Y_i^*)^2}{n(\sigma_C^2 + 1)}$ . Based on Duncan et al (2001; 2004)'s utility loss measure, we have the following:

$$UL_1 = \text{Var}(U_Y | \{y_i\}_{i=1}^n) = \frac{\sigma_C^2 \sum_{i=1}^n y_i^2}{n^2}$$

and

$$UL_2 = \text{Var}(U_{Y^2} | \{y_i\}_{i=1}^n) = \frac{[E(C^4) - (\sigma_C^2 + 1)^2] \sum_{i=1}^n y_i^4}{n^2(\sigma_C^2 + 1)^2}$$

The expression of  $UL_1$  means that, a set of noise candidates with equal variance will result in the same level of data utility loss for  $E(Y)$ . As a result, in this paper we use the following **overall utility loss measure**: When selecting among a set of noise candidates, we let the variances of these noise candidates to be equal. Then, we say a noise candidate will result in the lowest level of overall utility loss if it has the lowest  $UL_2$  value.

We note that the overall utility loss measure we use in this paper is simple and is for illustration purpose only. The primary purpose of this paper is the introduction of the proposed disclosure risk measure. In practice a data provider may use its own overall utility loss measure in conjunction with the disclosure risk measure we proposed for noise generating variable selection.

## 6 Simulations

In this section we present a simulation study. We show that our proposed disclosure risk measure could help a data provider to select an appropriate noise generating variable during data masking stage. We will use an R-U map (Duncan et al. 2001; 2004) to aid us with decision-makings in the simulation. We assume that the unbiased estimator and the correlation-attack estimator are the only sources of disclosure risk. We also comment on the protection levels offered by a few noise candidates which guarantee that  $R_{LW}$  is very small or is 0.

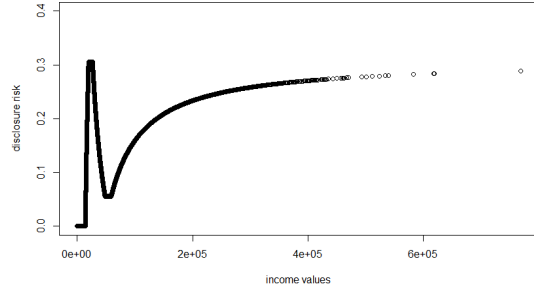
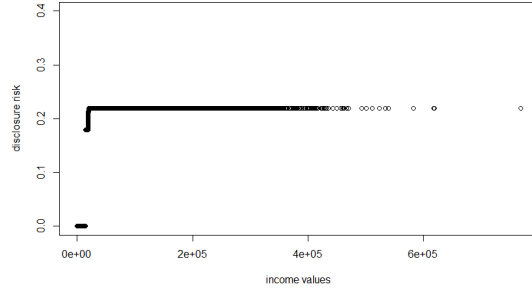
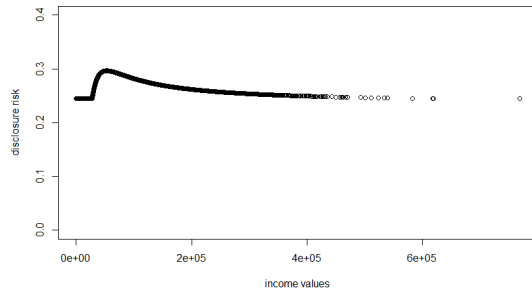
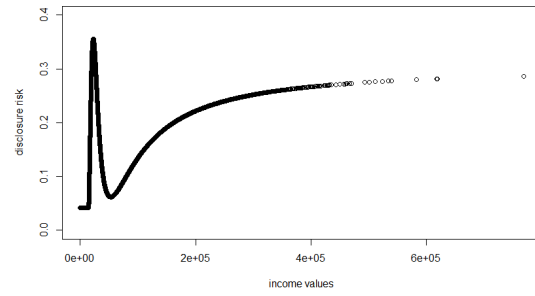
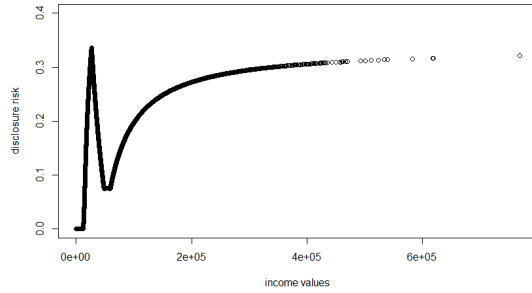
Following Klein et al. (2014), we use the public use data from the 2000 Current Population Survey (CPS) March supplement (available from <http://www.census.gov/cps/>). The entire data set contains household, family, and individual records. In this paper, we consider **positive household income values** under household income attribute as the original data. The original data contains 50661 positive observations ranging from 1 to 768742, with mean 53007 and variance 2411407246. The data is skewed to the right. In the following we denote the original data as  $\{y_i\}_{i=1}^{50661}$ .

We consider five noise candidates  $\{C_i\}_{i=1}^5$  with equal variance. We aim to use our proposed disclosure risk and overall utility loss measures to determine the best noise candidate for masking  $\{y_i\}_{i=1}^{50661}$ . We set the acceptance rule  $\delta = 0.1$  throughout this section. Among these noise candidates,  $C_1 \sim I_1 U_1 + (1 - I_1) U_2$ , where  $P(I_1 = 0) = P(I_1 = 1) = 0.5$ ,  $U_1 \sim U(0.5, 0.9)$  and  $U_2 \sim U(1.1, 1.5)$ ;  $C_2 \sim U(1 - 0.5\sqrt{93/75}, 1 + 0.5\sqrt{93/75})$ ;  $C_3$  follows a truncated normal distribution. That is, it follows  $N(1, 0.146135)$  which is truncated at 0.3 and 1.7, meaning that  $C_3$  has support  $(0.3, 1.7)$ ;  $C_4 \sim I_2 N_1 + (1 - I_2) N_2$ , where  $P(I_2 = 0) = P(I_2 = 1) = 0.5$ ,  $N_1 \sim N(0.7, 4/300)$  and  $N_2 \sim N(1.3, 4/300)$ ;  $C_5 \sim I_3 T_1 + (1 - I_3) T_2$ , where  $P(I_3 = 0) = P(I_3 = 1) = 0.5$ ,  $T_1$  and  $T_2$  are triangular random variables with three parameters  $(1.1 - \sqrt{\frac{9.6}{4}}, 0.9, 0.9)$  and  $(1.1, 1.1, 0.9 + \sqrt{\frac{9.6}{4}})$  respectively. The distributions of these noise candidates have been proposed or used in the literature for producing noise-multiplied data. Specifically,  $C_1$  follows a double uniform distribution, which was considered in the simulation study in Klein et al. (2014);  $C_2$  follows a uniform distribution, which was proposed and discussed in Sinha et al. (2011);  $C_3$  follows a truncated normal distribution, which was considered in Kim and Winkler (2003);  $C_4$  follows a bi-modal normal distribution, which was proposed in Lin and Wise (2012);  $C_5$  follows a truncated triangular distribution, which was proposed in Kim (2007) and Kim and Jeong (2008). Note that  $C_1$ ,  $C_4$  and  $C_5$  lead to  $R_{LW} = 0$  for all original observations. Therefore they seem to be good noise candidates if the unbiased estimator  $Y_i$  is the only source of disclosure risk for  $y_i$ . For  $C_4$ , there is a very small probability that it produces a negative noise. We simply ignore this fact in this simulation because it does not affect the simulation results.

Now we assume the role of the data provider and we aim to find an appropriate noise candidate to use during data masking stage. We assume the following **criteria for disclosure risk control**: a noise candidate needs to guarantee that  $\max(\{R(y_i, \delta)\}_{i=1}^n) < 0.3$ , i.e. the disclosure risk of any original value is less than 0.3. When multiple noise candidates satisfy the criteria for disclosure risk control, we say a noise candidate is better than the others if it offers the lowest average disclosure risk  $\{R(y_i, \delta)\}_{i=1}^n$ . We use **Approach 1** described in Section 5 to estimate  $R_\rho$  using  $R_{cor}$ . We also use the **overall utility loss measure** introduced in Section 5 to compare the levels of overall utility loss of the noise candidates. That is, since all the noise candidates have the same variance, a noise candidate with a lower  $UL_2$  has a lower level of overall utility loss under our overall utility loss measure.

The disclosure risk plots of  $\{R(y_i, 0.1)\}_{i=1}^{50661}$  for all five noise candidates are given in Figure 1. To comment on the plots, we see that different noise candidates protect original values differently. For instance, we see that  $C_2$  (Figure 1(b)) offers uniform protections to most observations while the others do not. **The noise candidates  $C_1$ ,  $C_5$ , which were previously thought to be good noise candidates because they guarantee that  $R_{LW} = 0$ , are actually not that ideal when our proposed disclosure risk measure  $R(y_i, \delta)$  is used.** For these noise candidates,  $R(y_i, \delta) = R_{cor}(y_i, \delta)$ , i.e. the disclosure risk comes entirely from the correlation-attack. We see that in this simulation, observations around the sample mean or extremely large are not protected well by these noise candidates. It may suggest that due to the correlation-attack, the merit of those noise candidates with  $R_{LW} = 0$  may need to be reconsidered even if the data intruder may only rely on two pieces of information  $y^*$  (noise-multiplied data) and  $\sigma_C^2$  (variance of noise) to attack an original value.

For each noise candidate, we also considered the average disclosure risk  $\overline{\{R(y_i, \delta)\}_{i=1}^{50661}}$  and the overall level of data utility loss it produces. We use an R-U map to visualize the utility-risk tradeoffs, which is presented in Figure 2. In the plot,  $C_1$  and  $C_4$  overlaps, meaning that they provide similar utility-risk tradeoffs. We see that, the overall data utility loss levels of these noise candidates are very similar. The average disclosure risks are very low for  $C_1$ ,  $C_4$  and  $C_5$ , followed by  $C_2$  and  $C_3$ .

(a) disclosure risk plot for  $C_1$ .(b) disclosure risk plot for  $C_2$ .(c) disclosure risk plot for  $C_3$ .(d) disclosure risk plot for  $C_4$ .(e) disclosure risk plot for  $C_5$ .Fig. 1: Disclosure risks  $\{R(y_i, 0.1)\}_{i=1}^{50661}$  plots for each noise candidate.

Based on all these results, we could make a decision about which noise candidate is appropriate to use in this context. Based on Figure 1, we see that  $C_1$ ,  $C_4$  and  $C_5$  do not satisfy our criteria for disclosure risk control because they cannot guarantee that every original observation has a disclosure risk less than 0.3. Therefore, we could only choose between  $C_2$  and  $C_3$ . From Figure 2, we see that  $C_2$  results in a much lower average disclosure risk level than  $C_3$  at the expense of only a slightly higher level of overall data utility loss. As a result,  $C_2$  seems to be the best choice in this context.

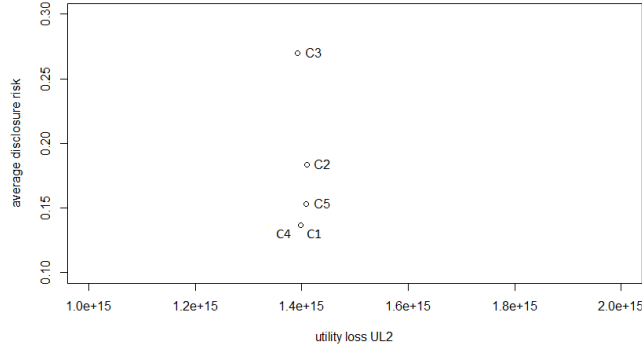


Fig. 2: Utility-Risk tradeoffs of the noise candidates.

## 7 Conclusion and future work

In this paper we showed the correlation-attack which could be used by data intruders to attack noise-multiplied data. The correlation-attack only uses information of noise-multiplied data itself and the variance of noise terms, hence could be applied to attack noise-multiplied data in most situations. Comparison between the correlation-attack estimator and the unbiased estimator which is another common attacking estimator were made and the result showed that the disclosure risks from both estimators need to be controlled. Correspondingly, we proposed a disclosure risk measure for data providers to evaluate the disclosure risks against both attacking estimators simultaneously for each noise candidate. The proposed disclosure risk measure could be used with other overall data utility loss measures to help data providers with decision-makings on noise candidates selection during data masking stage.

For an original value, a noise candidate which results in a low disclosure risk from the unbiased estimator may not result in a low disclosure risk from the correlation-attack estimator and vice versa. Similarly, for a set of original data, a noise candidate which produces a low average level of disclosure risk from the unbiased estimator may produce a high average level of disclosure risk from the correlation-attack estimator and vice versa. For instance, if the original data follows  $LN(5, 0.12^2)$ , we can show that  $C_4$  results in a low average level of disclosure risk from the unbiased estimator (0.0414) but a high average level of disclosure risk from the correlation-attack estimator (roughly 0.619). Similarly, we can show that the noise candidate  $C_3$  will result in a very low average level of disclosure risk from the correlation-attack estimator (actually 0) but a high average level of disclosure risk from the unbiased estimator (0.221). An ideal noise candidate should protect against both estimators effectively. It might be interesting to find out these noise candidates given different distributions of original data in the future.

Identifications of other attacking strategies are necessary for disclosure risk control of noise-multiplied microdata. For instance, it may be the case that in a microdata, multiple attributes are protected by multiplicative noises. This masking scenario is considered in Nayak et al. (2011) and Lin and Wise (2012). In that case, the idea of the regression-attack proposed in Li and Sarkar (2011) might be used to attack a noise-multiplied value. That is, to attack  $y_{ij}$ , which is the original value of the  $i$ -th record on the  $j$ -th attribute, a data intruder may regress  $y_{ij}$  on other noise-multiplied values  $(y_{i1}^*, y_{i2}^*, \dots, y_{ip}^*)$ , where  $p$  is the number of regressors used to attack  $y_{ij}$ . The idea could be explored in the future.



**Acknowledgements** This research has been conducted with the support of the Australian Government Research Training Program Scholarship.

## References

1. Agrawal, R. and Aggarwal, C. (2001) On the design and quantification of privacy preserving data mining algorithms. In *Proceedings of the 20th Symposium on Principles of Database Systems*, Santa Barbara, California, USA.
2. Agrawal, R. and Srikant, R., Privacy preserving data mining. In *Proceedings of the ACM SIGMOD*, pp. 439-450 (2000)
3. Brand, R., Microdata protection through noise addition. In *Inference Control in Statistical Databases*, vol. 2316 of LNCS. Springer Berlin Heidelberg, pp. 61-74 (2002)
4. Domingo-Ferrer, J., Seb , F., and Castell -Roca, J., On the security of noise addition for privacy in statistical databases. *Lecture notes in computer science*, **3050**, 149-161 (2004)
5. Domingo-Ferrer, J. and Torra, V., Disclosure Protection Methods and Information Loss for Microdata. In *Confidentiality, Disclosure and Data Access: Theory and Practical Applications for Statistical Agencies* (eds. Doyle P., Lane J.I., Theeuwes J.J.M. and Zayatz L.), pp.91-110 (2001)
6. Duncan, G., Keller-McNulty, S., and Stokes, S., Disclosure risk vs. data utility: the R-U confidentiality map. *Technical Report LA-UR-01-6428*, Los Alamos National Laboratory, Statistical Sciences Group, Los Alamos, New Mexico (2001)
7. Duncan, G., Keller-McNulty, S., and Stokes, S., Database security and confidentiality: examining disclosure risk vs. data utility through the R-U confidentiality map. *Technical Report Number 142*, National Institute of Statistical Science (2004)
8. Evans, T., Effects on trend statistics of the use of multiplicative noise for disclosure limitation. U.S. Bureau of the Census, <http://www.census.gov/srd/sdc/papers.html> (1996)
9. Hwang, J. T., Multiplicative errors-in-variables models with applications to recent data released by the U.S. Department of Energy. *Journal of American Statistical Association*, **81**, 680-688 (1986)
10. Kargupta, H., Datta, S., Wang, Q., and Sivakumar, K., On the privacy preserving properties of random data perturbation techniques. *Proceedings of the 3rd IEEE international conference on data mining*, Melbourne, 99-106 (2003)
11. Klein, M., Mathew, T., and Sinha, B., Noise multiplication for statistical disclosure control of extreme values in log-normal regression samples. *Journal of Privacy and Confidentiality*, **6**, 77-125 (2014)
12. Kim, J. J., Application of truncated triangular and trapezoidal distributions for developing multiplicative noise. Proceedings of the Survey Methods Research Section, *American Statistical Association*, CD Rom (2007)
13. Kim, J., Jeong, D. M., Truncated triangular distribution for multiplicative noise and domain estimation. Section on Government Statistics-JSM 2008, 1023-1030 (2008)
14. Kim, J.J., Winkler, W. E., Masking microdata files. *American Statistical Association, Proceedings of the Section on Survey Research Methods*, 114-119 (1995)
15. Kim, J.J., Winkler, W. E., Multiplicative noise for masking continuous data. Statistical Research Division, Research Report Series(Statistics #2003-01). U.S. Census Bureau (2003)
16. Li, X. B. and Sarkar, S., Protecting Privacy Against Regression Attacks in Predictive Data Mining. International Conference on Information Systems, Icis 2011, Shanghai, China (2011)
17. Li, X. B. and Sarkar, S., Class-restricted clustering and microperturbation for data privacy. *Management Science*, 59(4), 796-812 (2013)
18. Lin, Y. X. and Wise, P., Estimation of regression parameters from noise multiplied data. *Journal of Privacy and Confidentiality*, **4**, 61-94 (2012)
19. Liu, K., Giannella, C., and Kargupta, H., A Survey of Attack Techniques on Privacy-Preserving Data Perturbation Methods. *Privacy-Preserving Data Mining, vol.34 of the series Advances in Database Systems*, 359-381 (2008)
20. Ma, Y., Lin, Y. X., Chipperfield, J. O., Newman, J., and Leaver, V., A new algorithm for protecting aggregated business microdata via a remote system. *Privacy in Statistical Databases*, LNCS, 9867, 210-221 (2016)
21. Muralidhar, K. and Domingo-Ferrer, J., Rank-Based record linkage for re-identification risk assessment. *Privacy in Statistical Databases*, LNCS, 9867, 225-236 (2016)
22. Nayak, T. K., Sinha, B., and Zayatz, L., Statistical properties of multiplicative noise masking for confidentiality protection. *Journal of Official Statistics*, **27**, No.3, 527-544 (2011)
23. Oganyan, A. and Karr, A., Masking methods that preserve positivity constraints in microdata. *Journal of Statistical Planning and Inference*, **141**, 31-41 (2011)
24. Rubin, D. B., Discussion of statistical disclosure limitation. *Journal of Official Statistics*, **9**, 461-468 (1993)
25. Sinha, B., Nayak, T.K., and Zayatz, L., Privacy protection and quantile estimation from noise multiplied data. *Sankhya B*, **73**, No. 2, 297-315 (2011)
26. Yancey, W.E., Winkler, W.E., and Creedy, R.H., Disclosure risk assessment in perturbative micro-data protection. *Inference Control in Statistical Databases* (ed. J. Domingo-Ferrer), New York: Springer, 135-151 (2002)